



THE UNIVERSITY of EDINBURGH

## Edinburgh Research Explorer

### Meta-analysis of up to 622,409 individuals identifies 40 novel smoking behaviour associated genetic loci

#### Citation for published version:

Erzurumluoglu, AM, Liu, M, Jackson, VE, Barnes, DR, Datta, G, Melbourne, CA, Young, RD, Batini, C, Surendran, P, Jiang, T, Adnan, SD, Afaq, S, Agrawal, A, Altmaier, E, Antoniou, AC, Asselbergs, FW, Baumbach, C, Bierut, L, Bertelsen, S, Boehnke, M, Bots, ML, Brazel, DM, Chambers, JC, Chang-claude, J, Chen, C, Corley, J, Chou, Y, David, SP, De Boer, RA, De Leeuw, CA, Dennis, JG, Dominiczak, AF, Dunning, AM, Easton, DF, Eaton, C, Elliott, P, Evangelou, E, Faul, JD, Foroud, T, Goate, A, Gong, J, Grabe, HJ, Haessler, J, Haiman, C, Hallmans, G, Hammerschlag, AR, Harris, SE, Hattersley, A, Heath, A, Hsu, C, Iacono, WG, Kanoni, S, Kapoor, M, Kaprio, J, Kardia, SL, Karpe, F, Kontto, J, Kooner, JS, Kooperberg, C, Kuulasmaa, K, Laakso, M, Lai, D, Langenberg, C, Le, N, Lettre, G, Loukola, A, Luan, J, Madden, PAF, Mangino, M, Marioni, RE, Marouli, E, Marten, J, Martin, NG, McGue, M, Michailidou, K, Mihailov, E, Moayyeri, A, Moitry, M, Müller-nurasyid, M, Naheed, A, Nauck, M, Neville, MJ, Nielsen, SF, North, K, Perola, M, Pharoah, PDP, Pistis, G, Polderman, TJ, Posthuma, D, Poulter, N, Qaiser, B, Rasheed, A, Reiner, A, Renström, F, Rice, J, Rohde, R, Rolandsson, O, Samani, NJ, Samuel, M, Schlessinger, D, Scholte, SH, Scott, RA, Sever, P, Shao, Y, Shrine, N, Smith, JA, Starr, JM, Stirrups, K, Stram, D, Stringham, HM, Tachmazidou, I, Tardif, J, Thompson, DJ, Tindle, HA, Tragante, V, Trompet, S, Turcot, V, Tyrrell, J, Vaartjes, I, Van Der Leij, AR, Van Der Meer, P, Varga, TV, Verweij, N, Völzke, H, Wareham, NJ, Warren, HR, Weir, DR, Weiss, S, Wetherill, L, Yaghootkar, H, Yavas, E, Jiang, Y, Chen, F, Zhan, X, Zhang, W, Zhao, W, Zhao, W, Zhou, K, Amouyel, P, Blankenberg, S, Caulfield, MJ, Chowdhury, R, Cucca, F, Deary, IJ, Deloukas, P, Di Angelantonio, E, Ferrario, M, Ferrières, J, Franks, PW, Frayling, TM, Frossard, P, Hall, IP, Hayward, C, Jansson, J, Jukema, JW, Kee, F, Männistö, S, Metspalu, A, Munroe, PB, Nordestgaard, BG, Palmer, CNA, Salomaa, V, Sattar, N, Spector, T, Strachan, DP, Van Der Harst, P, Zeggini, E, Saleheen, D, Butterworth, AS, Wain, LV, Abecasis, GR, Danesh, J, Tobin, MD, Vrieze, S, Liu, DJ & Howson, JMM 2019, 'Meta-analysis of up to 622,409 individuals identifies 40 novel smoking behaviour associated genetic loci', *Molecular Psychiatry*. <https://doi.org/10.1038/s41380-018-0313-0>

#### Digital Object Identifier (DOI):

[10.1038/s41380-018-0313-0](https://doi.org/10.1038/s41380-018-0313-0)

#### Link:

[Link to publication record in Edinburgh Research Explorer](#)

#### Document Version:

Publisher's PDF, also known as Version of record

#### Published In:

Molecular Psychiatry

#### General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.





# Meta-analysis of up to 622,409 individuals identifies 40 novel smoking behaviour associated genetic loci

A. Mesut Erzurumluoglu<sup>1</sup> et al.

Received: 2 May 2018 / Revised: 30 September 2018 / Accepted: 14 November 2018  
© The Author(s) 2019. This article is published with open access

## Abstract

Smoking is a major heritable and modifiable risk factor for many diseases, including cancer, common respiratory disorders and cardiovascular diseases. Fourteen genetic loci have previously been associated with smoking behaviour-related traits. We tested up to 235,116 single nucleotide variants (SNVs) on the exome-array for association with smoking initiation, cigarettes per day, pack-years, and smoking cessation in a fixed effects meta-analysis of up to 61 studies (up to 346,813 participants). In a subset of 112,811 participants, a further one million SNVs were also genotyped and tested for association with the four smoking behaviour traits. SNV-trait associations with  $P < 5 \times 10^{-8}$  in either analysis were taken forward for replication in up to 275,596 independent participants from UK Biobank. Lastly, a meta-analysis of the discovery and replication studies was performed. Sixteen SNVs were associated with at least one of the smoking behaviour traits ( $P < 5 \times 10^{-8}$ ) in the discovery samples. Ten novel SNVs, including rs12616219 near *TMEM182*, were followed-up and five of them (rs462779 in *REV3L*, rs12780116 in *CNNM2*, rs1190736 in *GPR101*, rs11539157 in *PJA1*, and rs12616219 near *TMEM182*) replicated at a Bonferroni significance threshold ( $P < 4.5 \times 10^{-3}$ ) with consistent direction of effect. A further 35 SNVs were associated with smoking behaviour traits in the discovery plus replication meta-analysis (up to 622,409 participants) including a rare SNV, rs150493199, in *CCDC141* and two low-frequency SNVs in *CEP350* and *HDGFRP2*. Functional follow-up implied that decreased expression of *REV3L* may lower the probability of smoking initiation. The novel loci will facilitate understanding the genetic aetiology of smoking behaviour and may lead to the identification of potential drug targets for smoking prevention and/or cessation.

## Introduction

Smoking is a major risk factor for many diseases, including common respiratory disorders such as chronic obstructive pulmonary disease (COPD) [1, 2], cancer [3] and

cardiovascular diseases [4], and is reported to cause 1 in 10 premature deaths worldwide [5]. A greater understanding of the genetic aetiology of smoking behaviour has the potential to lead to new therapeutic interventions to aid smoking prevention and cessation, and thereby reduce the global burden of such diseases.

Previous genome-wide association studies (GWASs) identified 14 common SNVs [1, 6–12] (with minor allele frequency, MAF > 0.01) robustly associated with smoking behaviour-related traits ( $P < 5 \times 10^{-8}$ ). The 15q25 (*CHRNA3/5-CHRNA4*) region has the largest effect, explaining ~1% and 4–5% of the phenotypic variance of smoking quantity [13] and cotinine, a biomarker of nicotine intake [14], respectively. Overall, genetic loci identified to date explain ~2% of the estimated genetic heritability of smoking behaviour [6], which is reported to be between 40–60% [15–17]. A recent study suggested that an important proportion (~3.3%) of the phenotypic variance of smoking behaviour-related traits was explained by rare nonsynonymous variants (MAF < 0.01)

These authors contributed equally and share the first author position:  
A. Mesut Erzurumluoglu, Mengzhen Liu, Victoria E. Jackson

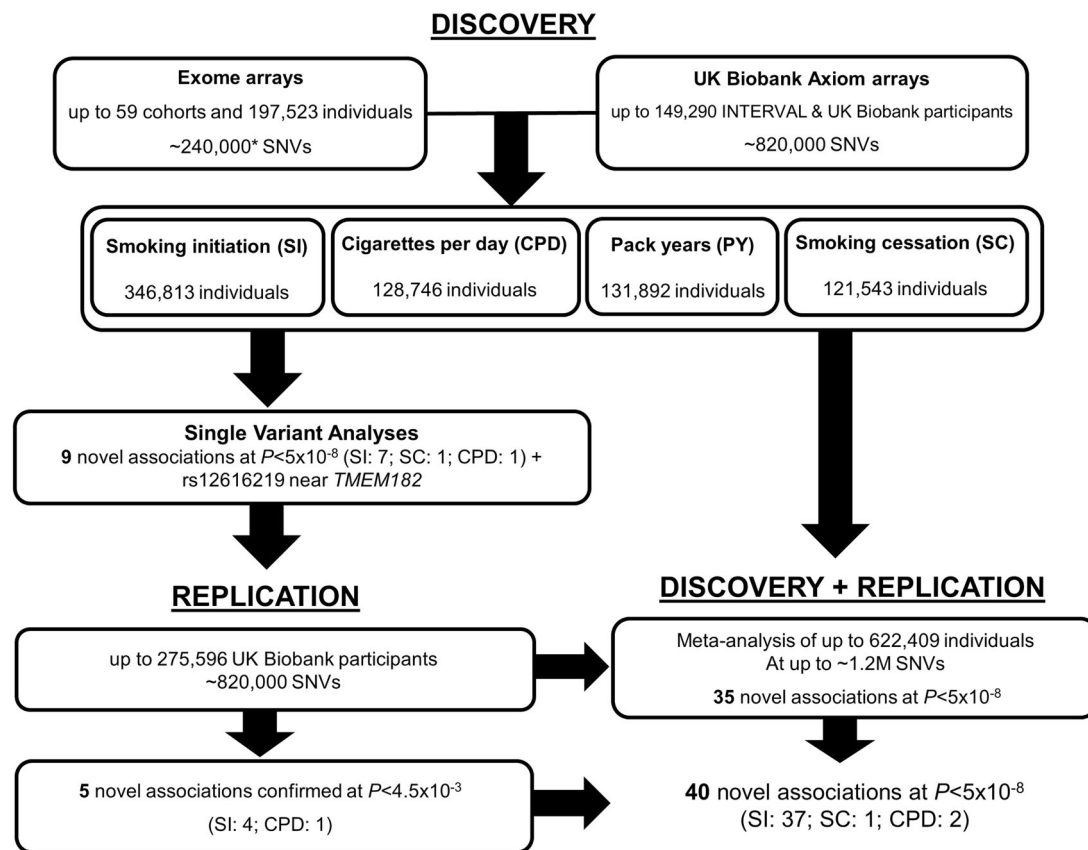
These authors contributed equally and share the last author position:  
Martin D. Tobin, Scott Vrieze, Dajiang J. Liu, Joanna M. M. Howson

**Supplementary information** The online version of this article (<https://doi.org/10.1038/s41380-018-0313-0>) contains supplementary material, which is available to authorized users.

✉ Dajiang J. Liu  
dxl46@psu.edu

✉ Joanna M. M. Howson  
jmmh2@medschl.cam.ac.uk

Extended author information available on the last page of the article.



**Fig. 1** Study design including the discovery and replication stages. NB: Gene-based studies, conditional analyses, and replication in African American ancestry samples not shown here for clarity. \*GFG and NAGOZALC studies contributed additional custom content

[18]. Hence, well-powered studies of rare variants are needed.

To investigate the effect of rare coding variants on smoking behaviour, we studied 346,813 participants (of which 324,851 were of European ancestry) from 61 cohorts (Supp. Tables 1 and 2) at up to 235,116 SNVs from the exome array. As we had access to UK Biobank, we also interrogated SNVs present on the UK Biobank and UK BiLEVE Axiom arrays to identify additional associations across the genome beyond the exome array. To our knowledge, these datasets are an order of magnitude larger than the previous studies [6], and constitute the most powerful exome-array study of smoking behaviour to date.

## Materials and methods

### Participants

Our study combined study-level summary association data from up to 59 studies of European ancestry and two studies of South Asian ancestry from three consortia (Consortium for Genetics of Smoking Behaviour (CGSB), GWAS & Sequencing Consortium of Alcohol and Nicotine use

(GSCAN) and the Coronary Heart Disease (CHD) Exome+consortium), INTERVAL and UK Biobank. In total, up to 324,851 individuals of European ancestry and 21,962 South Asian individuals were analysed in the discovery stage (Fig. 1). Further information about the participating cohorts and consortia is given in Supp. Table 1 and the Supp. Material. All participants provided written informed consent and studies were approved by local Research Ethics Committees and/or Institutional Review boards.

### Phenotypes

We chose to analyse the following four smoking behaviour-related traits because of their broad availability in existing epidemiological and medical studies, as well as their biological relevance for addiction behaviours:

- Smoking initiation (binary trait: ever vs never smokers). Ever smokers were defined as individuals who have smoked >99 cigarettes in their lifetime, which is consistent with the definition by the Centre for Disease Control [19];
- Cigarettes per day (CPD; quantitative trait: average number of cigarettes smoked per day by ever smokers);

- iii. Pack-years (quantitative trait; Packs per day x Years smoked, with a pack defined as 20 cigarettes); years smoked is typically formed from age at smoking commencement to current age for current smokers or age at cessation for former smokers.
- iv. Smoking cessation (binary trait: former vs current smokers).

In UK Biobank, phenotypes were defined using phenotype codes 1239, 1249, and 2644 for smoking initiation and smoking cessation, and 1239, 3436, 3456 for CPD and pack-years. CPD was inverse normal transformed in the CHD Exome+, INTERVAL and CGSB studies and categorised (1–10, 11–20, 21–30, and 31+ CPD) by the GSCAN studies and UK Biobank (Supp. Table 2). All studies performed an inverse normal transformation of pack-years. Summary statistics of study level phenotype distributions are provided in Supp. Table 1.

### Genotyping and quality control

Fifty-nine cohorts were genotyped using exome arrays (up to 235,116 SNVs) and two (UK Biobank and INTERVAL) were genotyped using Axiom Biobank Arrays (up to 820,000 SNVs; Supp. Table 2). In total, ~1.06M SNVs were analysed including ~64,000 SNVs on both the Axiom and Exome Arrays. Furthermore, two studies (NAGOZALC and GFG) genotyped their participants using arrays with custom content, increasing the total number of variants analysed to 1,207,583 SNVs. Individual studies performed quality control (QC; Supp. Material, Supp. Table 2) and additional QC was conducted centrally (i) to ensure alleles were consistently aligned, (ii) that there were no major sample overlaps between contributing studies, and (iii) variants conformed to Hardy–Weinberg equilibrium and call rate thresholds. We also examined the distribution of the effect sizes and test statistics across cohorts to ensure the test statistics were well-calibrated.

### Study level analyses

Each study (including the case-cohort studies [20]) undertook analyses of up to four smoking traits using RAREMETALWORKER [21] or RVTESTS [22] (Supp. Table 2), which generated single variant score statistics and their covariance matrices within sliding windows of 1Mb. CPD and pack-years were analysed using linear models or linear mixed models. Smoking initiation and smoking cessation were analysed using logistic models or linear mixed models. All studies adjusted each trait for age, sex, at least three genetic principal components and any study-specific covariates (Supp. Table 2). Chromosome X variants were analysed using the above-described approach, but coding

males as 0/2. This coding scheme ensures that on average females and males have equal dosages and so is optimal for genes that are inactivated (due to X chromosome inactivation) and is valid for genes that do not undergo X chromosome activation. Males and females were analysed together adjusting for sex as a covariate.

### Single variant meta-analyses

Fixed effects meta-analyses across the individual contributing studies of single variant associations were undertaken using the Cochran-Mantel-Haenszel method in RAREMETAL. Z-score statistics were used in the meta-analysis to ensure that the association results are robust against potentially different units of measurement in the phenotype definitions across studies [23]. We performed genomic control correction on the meta-analysis results. Variants with  $P < 1 \times 10^{-6}$  in tests of heterogeneity were excluded. Variants with  $P \leq 5 \times 10^{-8}$  were taken forward for replication. In addition, rs12616219 was also taken forward for replication as its  $P$ -value was very close to this threshold (smoking initiation,  $P = 5.49 \times 10^{-8}$ ). None of the rare SNVs were genome-wide significant, therefore we also took forward the rare variant with the smallest association  $P$ -value, rs141611945 ( $P = 2.95 \times 10^{-7}$ ; MAF < 0.0001).

### Replication and combined meta-analysis of discovery and replication data

As UK biobank genetic data were released in two phases, we took the opportunity to replicate findings from the discovery stage in a further 275,596 individuals made available in the phase two release of UK Biobank genetic data. To avoid potential relatedness between discovery and replication samples, the replication samples were screened and individuals with relatedness closer than second degree with the discovery sample in the UK Biobank were removed [24]. Phenotypes were defined in the same way as the discovery samples (described above). Since the exome array and the UK Biobank Axiom arrays do not fully overlap, we used both genotyped exome variants (approx. 64,000) as well as the additional ~90,000 well-imputed exome array variants from UK Biobank (imputation quality score >0.3) for replication of single variant and gene-based tests. The rare *ATF6* variant was absent from the UK Biobank array and is more prevalent in Africans (MAF = 0.01) than Europeans (MAF = 0.0007). Therefore, replication was sought in 1,437 individuals of African American-ancestry from the HRS and COGA studies. Analysis methods for replication cohorts were the same as for discovery cohorts, including methods to analyse chromosome X (Supp. Table 2). The criteria set for the replication were (i) the same direction of

effect as the discovery analysis and (ii)  $P \leq 0.0045$  in the replication studies (Bonferroni-adjusted for eleven SNVs at  $\alpha = 0.05$ ).

Finally, in order to fully utilise all available data, we carried out a combined meta-analysis of the discovery and replication samples across the exome array content using the same protocols mentioned above.

### Conditional analyses

To identify conditionally independent variants within previously reported and novel loci a sequential forward stepwise selection was performed [25]. A 1 MB region was defined around the reported or novel sentinel variant (500 kb either side) and conditional analyses performed with all variants within the region. If a conditionally independent variant was identified, ( $P < 5 \times 10^{-6}$ ; Bonferroni-adjusted for ~10,000 independent variants in the test region) the analysis was repeated conditioning on both the most significant conditionally independent variant and the sentinel variant. This stepwise approach was repeated (conditioning on the variants identified in current and earlier iterations) until there were no variants remaining in the region that were conditionally independent. The same protocol was followed for the novel SNVs identified in this study.

### Gene-based analyses

For discovery gene-based meta-analyses, we utilised three statistical methods as part of the RAREMETAL package: the Weighted Sum Test (WST) [26], the burden test [27] and the Sequence Kernel Association test (SKAT) [28]. EPACTS (v.3.3.0) [29] was used to annotate variants (for use in gene-based meta-analyses), as recommended by RAREMETAL. Two MAF cut-offs were used, one used low-frequency (MAF < 0.05) and rare variants, the second only used rare variants (MAF < 0.01). Nonsynonymous, stop gain, splice site, start gain, start loss, stop loss, and synonymous variants were selected for inclusion. A sensitivity analysis to exclusion of synonymous variants was also performed. Gene-level associations with  $P < 8 \times 10^{-7}$  were deemed statistically significant (Bonferroni-adjusted for ~20,000 genes and three tests at  $\alpha = 0.05$ ). To examine if the gene associations were driven by a single variant, the gene tests were conducted conditional on the SNV with the smallest  $P$ -value in the gene, using the shared single variant association statistic and covariance matrices [21, 25].

### Mendelian randomization analyses

To evaluate the causal effect of SI and CPD on BMI, schizophrenia and educational attainment (EA), we conducted Mendelian randomization (MR) analyses using three

complementary approaches available in MR-Base [30]: inverse variance weighted regression [31], MR-Egger [32, 33], and weighted median [34]. We used both the previously reported smoking-associated SNVs and the SNVs from the current report (as provided in Tables 1–3 and Supp. Table 3) as instrumental variables. The BMI [35], schizophrenia [36] and educational attainment [37] data came from previously published publicly available data. To assess possible reverse causation, we also used outcome associated SNVs as instrumental variables and conducted MR analyses using SI and CPD as outcome. We considered  $P < 0.05/3 = 0.017$  as statistically significant (Bonferroni-adjusted for three traits).

### In silico functional follow up of associated SNVs

To identify whether the (replicated) SNVs identified here affected other traits, we queried the GWAS Catalog [38] (version: e91/28/02/2018, downloaded on 01/03/18) for genome-wide significant ( $P < 5 \times 10^{-8}$ ) associations using all proxy SNVs ( $r^2 \geq 0.8$ ) within 2 Mb of the top variant in our study.

eQTL lookups were carried out in the 13 brain tissues available in GTEx V7 [39], Brain xQTL (dorsolateral prefrontal cortex) [40] and BRAINEAC [41] databases, all of which had undergone QC by the individual studies. We did not perform additional QC on these data. In brief, GTEx used Storey's  $q$ -value method to correct the FDR for testing multiple transcripts based upon the empirical  $P$ -values for the most significant SNV for each transcript [43, 42]. BRAINEAC calculated the number of tests per transcript and used Benjamini–Hochberg procedure to calculate FDR per transcript using a FDR < 1% as significant. BRAINxQTL used  $P < 8 \times 10^{-8}$  as a cut-off for significance for any given transcript. SNVs that met the study specific significance and FDR thresholds, which were in LD ( $r^2 > 0.8$  in 1000 Genomes Europeans) with the top eQTL or the sentinel eQTL for a given tissue/transcript combination were considered significant. The genes implicated by these eQTL databases and/or coding changes (e.g., missense and nonsense SNVs) were put into ConsensusPathDB [44] to identify whether these genes were over-represented in any known biological pathways. Replicated missense SNVs were also put into PolyPhen-2 [45] and FATHMM (unweighted) [46] to obtain variant effect prediction.

## Results

### Single variant associations

In the discovery meta-analyses, we identified 15 common SNVs that were genome-wide significant ( $P < 5 \times 10^{-8}$ ) for one or more of the smoking behaviour traits, of which 9



**Table 1** Association results for SNVs identified in single variant association meta-analyses and taken forward to replication are provided

dbSNP ID (Exome ID)	Chr:Pos	EA/OA	Gene	Consequence	Trait	Discovery stage			Replication stage	
						N	EAF	DoE P-value	Beta (SE)	P-value
rs141611945 (exm118559)	1:161771868	G/A	ATF6	Missense	CPD	128,746	0.0065% MAC = 9	+ 2.95 × 10 <sup>-7</sup>	0.184 (0.169)	*P = 0.276 in African American samples
rs1190736 ** (exm1659559)	X:136113464	A/C	GPR101	Missense	CPD (PY)	99,037 (96,824)	46.6% (47.0%)	- 1.40 × 10 <sup>-11</sup> (4.98 × 10 <sup>-9</sup> )	-0.028 (0.0041) -0.027 (0.0049) -0.028 (0.0073)	All samples: <b>8.20 × 10<sup>-12</sup></b> (2.70 × 10 <sup>-11</sup> ) Males only: <b>1.90 × 10<sup>-8</sup></b> (6.0 × 10 <sup>-8</sup> ) Females only: <b>1.10 × 10<sup>-4</sup></b> (7.1 × 10 <sup>-4</sup> )
rs462779 (exm572256)	6:111695887	A/G	REV3L	Missense	SI	346,682	80.1%	- 4.52 × 10 <sup>-8</sup>	-0.023 (0.0034)	<b>9.7 × 10<sup>-12</sup></b>
rs216195 (exm1276230)	17:2203167	G/T	SMG6	Missense	SI	335,406	27.3%	- 2.80 × 10 <sup>-8</sup>	-0.008 (0.0029)	8.5 × 10 <sup>-3</sup>
rs11539157 (exm1643833)	X:68381264	A/C	P/JAI	Missense	SI	289,917	16.5%	+ 1.39 × 10 <sup>-11</sup>	0.022 (0.0026) 0.0158 (0.0033) 0.0185 (0.0039)	All samples: <b>5.40 × 10<sup>-17</sup></b> Males only: <b>1.30 × 10<sup>-6</sup></b> Females only: <b>2.20 × 10<sup>-6</sup></b>
<i>Non-Exome-chip SNVs</i>										
rs12616219	2:104352495	A/C	TMEM182	Intergenic	SI	112,811	46.4%	- 5.49 × 10 <sup>-8</sup>	-0.015 (0.0027)	<b>5.5 × 10<sup>-8</sup></b>
rs1150691	6:28168033	G/A	ZSCAN9	Missense	SI	112,811	34.8%	- 4.95 × 10 <sup>-8</sup>	-0.007 (0.0028)	8.0 × 10 <sup>-3</sup>
rs2841334	9:128122320	A/G	GAPVD1	Intronic	SI	112,811	20.9%	- 2.28 × 10 <sup>-8</sup>	-0.009 (0.0033)	7.5 × 10 <sup>-3</sup>
rs202664	22:41813886	C/T	TOB2	Intergenic	SC	51,043	19.9%	- 1.02 × 10 <sup>-8</sup>	-0.011 (0.0050)	2.1 × 10 <sup>-2</sup>
rs11895381	2:60053727	A/G	BCL11A	Intergenic	SI	112,811	34.2%	- 5.61 × 10 <sup>-9</sup>	-0.007 (0.0028)	1.2 × 10 <sup>-2</sup>
rs12780116	10:104821946	A/G	CNNM2	Intronic	SI	112,811	13.9%	+ 9.19 × 10 <sup>-10</sup>	0.017 (0.0039)	<b>1.1 × 10<sup>-5</sup></b>

Novel smoking trait associated SNVs that replicated with  $P < 0.005$  and had consistent direction of effect in discovery and replication are highlighted in bold. The replication sample size for smoking initiation (SI), cigarettes per day (CPD), pack-years (PY), and smoking cessation (SC) were 275,596, 80,015, 78,897, and 123,851 respectively. Chromosome (Chr) and position (Pos) for hg19 build 37. EA effect allele, OA other allele, Gene closest gene, N number of individuals, EAF effect allele frequency in the pooled samples, MAC minor allele count, DoE direction of effect, SE standard error. All SNVs had heterogeneity  $P > 0.02$  in the discovery stage. \*Replication was sought in 1,437 individuals of African American-ancestry from the HRS and COGA studies; \*\*The beta(se) for the association of rs1190736 with PY in the replication stage was -0.026 (0.0039)

**Table 2** Association results for novel SNVs identified in the combined meta-analysis of the discovery and replication cohorts

dbSNP ID (Exome-chip ID)	Chr:Pos	EA/OA	Gene	Consequence	Trait	EAF	Beta (se) in replication stage	P-value in combined meta-analysis (P-value in Discovery/Replication stage)	Notes
<i>Combining only genotyped Exome-chip content on the Axiom array</i>									
rs1514175	1:74991644	G/A	<i>TNNI3K</i>	Intronic	SI	0.57	−0.011 (0.003)	<b><math>5.42 \times 10^{-9}</math></b> ( $9.03 \times 10^{-5}$ ) $1.0 \times 10^{-5}$	Previously associated with BMI
rs7096169	10:104618695	G/A	<i>BORCS7</i> ( <i>CNNM2</i> <sup>#</sup> in Table 1)	Intronic	SI	0.31	0.016 (0.003)	<b><math>2.17 \times 10^{-13}</math></b> ( $3.38 \times 10^{-7}$ ) $7.3 \times 10^{-9}$	$r^2 = 0.28$ between rs7096169 and rs12780116 (Table 1) in 1000 Genomes EUR. Previously associated with Schizophrenia. rs7096169 an eQTL for <i>ARL3</i> , <i>BORCS7</i> , and <i>AS3MT</i> in ≥1 of the brain tissues in GTEx
rs2292239	12:56482180	G/T	<i>ERBB3</i>	Intronic	SI	0.66	0.0121 (0.003)	<b><math>2.78 \times 10^{-8}</math></b> ( $7.56 \times 10^{-5}$ ) $1.5 \times 10^{-5}$	Previously associated with type-1 diabetes and years of educational attainment. rs2292239 is an eQTL for <i>RPS26</i> and <i>SUOX</i> in ≥4 of the brain tissues in GTEx
rs216195	17:2203167	G/T	<i>SMG6</i> <sup>#</sup>	Missense	SI	0.29	−0.0076 (0.003)	<b><math>2.41 \times 10^{-9}</math></b> ( $2.80 \times 10^{-8}$ ) $8.5 \times 10^{-3}$	Same SNV as in Table 1
<i>Combining well-imputed Exome-chip content on the Axiom array</i>									
rs2960306 (exm383568)	4:2990499	T/G	<i>GRK4</i>	Missense	CPD	0.34	−0.024 (0.005)	<b><math>1.06 \times 10^{-9}</math></b> ( $3.99 \times 10^{-5}$ ) $3.8 \times 10^{-6}$	rs2960306 is an eQTL for <i>GRK4</i> in four of the brain tissues in GTEx
rs4908760	1:8526142	A/G	<i>REKE</i>	Intronic	SI	0.35	0.0078 (0.003)	<b><math>1.76 \times 10^{-8}</math></b> ( $3.36 \times 10^{-6}$ ) $4.7 \times 10^{-3}$	Previously associated with Vitiligo
rs692219 (exm127721)	1:179989584	C/G	<i>CEP350</i>	Missense	SI	0.028	−0.0257 (0.008)	<b><math>4.69 \times 10^{-9}</math></b> ( $1.08 \times 10^{-6}$ ) $1.3 \times 10^{-3}$	
rs11971186	7:126437897	G/A	<i>GRM8</i>	Intronic	SI	0.20	−0.0080 (0.003)	<b><math>1.45 \times 10^{-8}</math></b> ( $1.38 \times 10^{-6}$ ) $3.9 \times 10^{-3}$	
rs150493199 (exm249655)	2:179721072	A/T	<i>CCDC141</i>	Missense	SC	0.0098	0.048 (0.134)	<b><math>1.28 \times 10^{-8}</math></b> ( $6.45 \times 10^{-8}$ ) 0.72)	
<i>Non-Exome-chip SNVs</i>									
rs3001723	1:44037685	A/G	<i>PTPRF</i>	Intronic	SI	0.21	0.0159 (0.003)	<b><math>6.64 \times 10^{-11}</math></b> (0.00015/ $4.1 \times 10^{-8}$ )	Previously associated with Schizophrenia and Years of educational attainment
rs1937455	1:66416939	G/A	<i>PDE4B</i>	Intronic	SI	0.30	−0.0146 (0.0027)	<b><math>1.23 \times 10^{-9}</math></b> (0.00073/ $5.6 \times 10^{-8}$ )	
rs72720396	1:91191582	G/A	<i>BARHL2</i>	Intergenic	SI	0.16	−0.0150 (0.003)	<b><math>9.86 \times 10^{-9}</math></b> ( $5.63 \times 10^{-5}$ ) $1.9 \times 10^{-6}$	
rs6673752	1:154219177	C/G	<i>UBAP2L</i>	Intronic	SI	0.055	−0.027 (0.004)	<b><math>1.1 \times 10^{-11}</math></b> (NA/1.1 × $10^{-11}$ )	
rs2947411	2:614168	G/A	<i>TMEM18</i>	Intergenic	SI	0.83	0.0189 (0.004)	<b><math>4.97 \times 10^{-10}</math></b> (0.00017/ $7.1 \times 10^{-8}$ )	Previously associated with BMI
rs528301	2:45154908	A/G	<i>SIX3</i>	Intergenic	SI	0.38	0.0136 (0.002)	<b><math>4.12 \times 10^{-11}</math></b> ( $1.77 \times 10^{-6}$ ) $3.8 \times 10^{-7}$	
rs6738833	2:104150891	T/C	<i>TMEM182</i> <sup>#</sup>	Intergenic	SI	0.33	−0.018 (0.003)	<b><math>8.66 \times 10^{-14}</math></b> ( $1.63 \times 10^{-6}$ ) $4.4 \times 10^{-11}$ )	$r^2 = 0.69$ between rs6738833 and rs12616219 (Table 1) in European samples of the 1000 Genomes Project
rs13026471	2:137564022	T/C	<i>THSD7B</i>	Intronic	SI	0.18	0.0127 (0.003)	<b><math>2.45 \times 10^{-8}</math></b> (0.00028/ $3.0 \times 10^{-5}$ )	
rs6724928	2:156005991	C/T	<i>KCNJ3</i>	Intergenic	SI	0.32	−0.011 (0.003)	<b><math>4.47 \times 10^{-8}</math></b> (0.0019/4.8 × $10^{-5}$ )	
rs13022438	2:162800372	G/A	<i>SLC4A10</i>	Intronic	SI	0.27	0.0146 (0.003)	<b><math>1.41 \times 10^{-11}</math></b> (0.0005/ $8.1 \times 10^{-8}$ )	

Table 2 (continued)

dbSNP ID (Exome-chip ID)	Chr:Pos	EA/OA	Gene	Consequence	Trait	EA <sup>a</sup>	Beta (se) in replication stage	P-value in combined meta-analysis (P-value in Discovery/Replication stage)	Notes
rs1869244	3:5724531	A/G	<i>LOC105376939</i>	Intergenic	SI	0.32	0.0123 (0.003)	<b><math>2.76 \times 10^{-9}</math></b> (0.00040/ $4.1 \times 10^{-6}$ )	
rs35438712	3:85588205	T/C	<i>CADM2</i>	Intronic	SI	0.25	0.017 (0.003)	<b><math>1.99 \times 10^{-13}</math></b> ( $1.15 \times 10^{-5}$ / $3.2 \times 10^{-10}$ )	
rs6883351	5:22193967	T/C	<i>CDH12</i>	Intronic	SI	0.34	0.0129 (0.003)	<b><math>4.69 \times 10^{-8}</math></b> (0.0010/1.4 × $10^{-6}$ )	
rs6414946	5:87729711	C/A	<i>TMEM161B</i>	Intronic	SI	0.32	−0.0137 (0.003)	<b><math>5.27 \times 10^{-10}</math></b> ( $3.63 \times 10^{-5}$ / $2.8 \times 10^{-7}$ )	
rs11747772	5:166992708	C/T	<i>TENM2</i>	Intronic	SI	0.25	0.0144 (0.003)	<b><math>6.20 \times 10^{-9}</math></b> (0.011/2.2 × $10^{-7}$ )	
rs9320995	6:98726381	G/A	<i>POU3F2</i>	Intergenic	SI	0.18	0.0150 (0.003)	<b><math>1.70 \times 10^{-8}</math></b> (0.00079/ $6.1 \times 10^{-7}$ )	
rs10255516	7:1675621	G/A	<i>ELFN1</i>	Intergenic	SI	0.33	−0.0139 (0.003)	<b><math>2.86 \times 10^{-10}</math></b> (0.0021/ $1.8 \times 10^{-7}$ )	
rs10807839	7:3344629	G/A	<i>SDK1</i>	Intronic	SI	0.19	0.0162 (0.003)	<b><math>8.93 \times 10^{-11}</math></b> (0.0026/ $4.4 \times 10^{-8}$ )	
rs6965740	7:117514840	T/G	<i>CTTNBP2</i>	Intergenic	SI	0.31	−0.0126 (0.003)	<b><math>9.66 \times 10^{-9}</math></b> ( $5.56 \times 10^{-6}$ / $2.8 \times 10^{-6}$ )	
rs11776293	8:27418429	T/C	<i>EPHX2</i>	Intronic	SI	0.12	−0.0200 (0.003)	<b><math>2.23 \times 10^{-12}</math></b> (0.00011/ $8.9 \times 10^{-9}$ )	rs11776293 is an eQTL for <i>CHRNA2</i> in cerebellum in GTEx
rs1562612	8:59817068	G/A	<i>TOX</i>	Intronic	SI	0.35	−0.0112 (0.003)	<b><math>1.15 \times 10^{-9}</math></b> ( $1.42 \times 10^{-5}$ / $2.9 \times 10^{-5}$ )	
rs3857914	8:93184065	C/T	<i>RUNX1T1</i>	Intergenic	SI	0.19	0.0157 (0.003)	<b><math>1.54 \times 10^{-9}</math></b> (0.065/7.1 × $10^{-8}$ )	
rs2799849	9:86752641	C/T	<i>RMI1</i>	Intergenic	SI	0.22	−0.0156 (0.003)	<b><math>1.94 \times 10^{-8}</math></b> (0.026/4.8 × $10^{-8}$ )	
rs6482190	10:22037809	A/G	<i>LOC107984214</i>	Intronic	SI	0.17	0.0146 (0.003)	<b><math>8.85 \times 10^{-9}</math></b> (0.0021/9.5 × $10^{-7}$ )	
rs4523689	11:7950797	G/A	<i>OR10A6</i>	Intergenic	SI	0.27	−0.012 (0.003)	<b><math>7.77 \times 10^{-9}</math></b> (0.00030/ $2.2 \times 10^{-5}$ )	
rs933006	13:38350193	A/G	<i>TRPC4</i>	Intronic	SI	0.32	−0.0143 (0.003)	<b><math>3.50 \times 10^{-8}</math></b> (0.022/9.6 × $10^{-8}$ )	
rs557899	15:47643795	A/C	<i>SEMA6D</i>	Intronic	SI	0.26	0.0157 (0.003)	<b><math>2.99 \times 10^{-13}</math></b> ( $4.46 \times 10^{-5}$ / $1.0 \times 10^{-8}$ )	
rs76608582	19:4474725	A/C	<i>HDCFRP2</i>	Intronic	SI	0.029	−0.0360 (0.007)	<b><math>8.50 \times 10^{-9}</math></b> (0.012/4.3 × $10^{-8}$ )	

Chromosome (Chr) and position (Pos) for each SNV is given for hg19 build 37. Only SNVs reaching genome-wide significance ( $P < 5 \times 10^{-8}$ , in bold) in the combined meta-analysis are shown. Magnitude of the effect size estimates are not presented as traits were transformed in differently by the three consortia analysed. SNVs identified in the discovery stage of this study (see Table 1) are denoted #. The discovery sample size for smoking initiation (SI), CPD, pack-years (PY), and smoking cessation (SC) were 346,813, 128,746, 131,892, and 121,543, respectively; and the replication sample size for SI, CPD, PY, and SC were 275,596, 80,015, 78,897, and 123,851, respectively. NB: rs6673752 (intronic to *UBAP2L*) was not available in the discovery cohorts. EA effect allele, OA other allele, Beta(se) beta and standard error for association in the replication stage. All SNVs had heterogeneity  $P > 0.0001$

Bold font highlights the genome-wide significant P-values from the meta-analysis of discovery plus replication studies



**Table 3** Results from conditional analyses at previously reported smoking behaviour loci

Gene region	dbSNP ID	Chr:Pos	EA/OA	Consequence	Trait	EAF	P (unconditional)	SNV(s) conditioned on	Discovery Conditional P [DoE]	Conditional P in replication [DoE]
19q13 ( <i>RAB4B</i> )	rs8102683	19:41363765	C/T	Intergenic	CPD	74.8%	<b><math>4.53 \times 10^{-16}</math></b>	rs7937	<b><math>1.44 \times 10^{-13}</math></b> [±]	$3.5 \times 10^{-4}$ [±]
	rs28399442	19:41354458	A/C	Intronic ( <i>CYP2A6</i> )	CPD	1.3%	<b><math>2.27 \times 10^{-12}</math></b>	rs7937, rs8102683	<b><math>2.63 \times 10^{-12}</math></b> [±]	<b><math>8.1 \times 10^{-14}</math></b> [±]
	rs3865453	19:41338556	T/C	Intergenic	CPD	6.54%	<b><math>2.96 \times 10^{-12}</math></b>	rs7937, rs8102683, rs28399442	<b><math>4.96 \times 10^{-10}</math></b> [−]	<b><math>2.3 \times 10^{-13}</math></b> [−]
<i>TEX41-</i> <i>PABPC1P2</i>	rs11694518	2:146125523	T/C	Intergenic	SI	29.5%	<b><math>2.90 \times 10^{-9}</math></b>	rs101193706	$3.43 \times 10^{-7}$ [−]	<b><math>4.0 \times 10^{-31}</math></b> [−]
15q25 ( <i>CHRNA3</i> )	rs938682	15:78882925	A/G	Intronic ( <i>CHRNA3</i> )	CPD	76.4%	<b><math>1.83 \times 10^{-69}</math></b>	rs1051730	<b><math>7.77 \times 10^{-21}</math></b> [±]	<b><math>1.0 \times 10^{-13}</math></b> [±]

SNVs with  $P < 5 \times 10^{-8}$  are highlighted in bold. The discovery sample size for smoking initiation (SI) and CPD was 346,813 and 128,746, respectively. The replication sample size for SI and CPD were 275,596 and 80,015, respectively. *Chr* Chromosome, *Pos* position for hg19 build 37, *EA* effect allele, *OA* other allele, *EAF* effect allele frequency in the pooled samples, *DoE* Direction of effect

were novel (Table 1, Supp. Table 3). Seven novel loci were identified for smoking initiation, one for both CPD and pack-years and one for smoking cessation (Figs. 1, 2, Table 1 and Supp. Figure 1). Results for the significant loci were consistent across participating cohorts and there was at least nominal evidence of association ( $P < 0.05$ ) at the novel loci within each of the contributing consortia (Supp.-Table 4). Full association results for all novel SNVs across the four traits are provided in Supp. Table 5. No rare variants were genome-wide significant; the rare variant with the smallest  $P$ -value was a missense variant in *ATF6*, rs141611945 (MAF < 0.0001, CPD  $P = 2.95 \times 10^{-7}$ ).

Eleven SNVs (including rs12616219 near *TMEM182* with  $P = 5.49 \times 10^{-8}$ , and the rare variant, rs141611945) were taken forward for replication in independent samples (Table 1). The latest release of European UK Biobank individuals not included in the discovery stage (smoking initiation,  $n = 275,596$ ; smoking cessation  $n = 123,851$ ; CPD  $n = 80,015$ ; pack-years  $n = 78,897$ ), was used for replication of the common variants (Fig. 1). Five of the common variants replicated (four for smoking initiation and one with CPD and pack-years) at  $P < 0.0045$ . Two coding variants (rs11539157, rs1190736) were predicted to be ‘probably damaging’ by PolyPhen-2 and FATHMM. The remaining five SNVs were at least nominally associated ( $P < 0.01$ ) in the replication samples and had consistent direction of effect across discovery and replication. Replication for the rare variant rs141611945 could not be carried out in UK Biobank as the SNV nor its proxies ( $r^2 > 0.3$ ) were available. Thus we initiated replication in African American samples of the COGA ( $n = 476$ ) and HRS ( $n = 961$ ) cohorts (overall MAF ≈ 0.01). The direction of effect was consistent in the two replication cohorts and consistent with the discovery meta-analysis but a meta-analysis of the two replication cohorts yielded a  $P = 0.28$ . Further data are required to replicate this association.

We also performed a meta-analysis combining the discovery and replication samples (up to 622,409 individuals). LD score regression showed that the  $\lambda$  (intercept) for all traits was ~1.00, which indicated that confounding factors inflating the results was not an issue [47, 48]. The combined analysis identified 35 additional novel SNV-smoking trait associations, 33 with smoking initiation, one with CPD and one with smoking cessation at  $P < 5 \times 10^{-8}$  (Table 2). We note that among our four SNVs that did not replicate, rs216195 (in *SMG6*) was genome-wide significant in the combined meta-analysis of discovery and replication studies ( $P = 2.41 \times 10^{-9}$ ; Table 2).

We also calculated the phenotypic variance explained for novel and known variants. Results can be found in the ‘Calculation of Phenotypic Variance Explained’ section in the Supplementary Material.



**Fig. 2** A concentric Circos plot of the association results for smoking initiation (SI; outer ring), cigarettes per day (CPD) and smoking cessation (SC; inner ring) for chromosomes 1–22 (Pack-years results, which can be found in Supp. Figure 1, are omitted for clarity). Each dot represents a SNV, with the X and Y axes corresponding to genomic location in Mb and  $-\log_{10}P$ -values, respectively. Labels show the nearest gene to the novel sentinel variants identified in the discovery

stage and taken forward to replication. The top signals were truncated at  $10^{-10}$  for clarity. Novel and previously reported signals are highlighted in red and dark blue, respectively. Grey rings on the y-axis increase by increments of 2 (initial ring corresponding to  $P = 0.001$ , then 0.00001 etc.); and the outer and inner red rings correspond to the genome-wide significance level ( $P = 5 \times 10^{-8}$ ) and  $P = 5 \times 10^{-7}$ , respectively. Image was created using Circos (v0.65)

### Associations at known smoking behaviour loci

We assessed evidence for associations at the 14 SNVs previously reported for smoking behaviour-related traits. Seven were genotyped on the exome array and proxies ( $r^2 > 0.3$ ;  $\pm 2$  Mb) were identified for the remaining seven (Supp. Table 3). All showed nominal evidence of association at

$P < 0.05$  and six of these were genome-wide significant in the meta-analysis of the trait for which it was previously reported (Supp. Tables 3 and 5).

Conditional analyses identified five independent associations within three previously reported loci and all five replicated (Table 3). At the 19q13 (*RAB4B*) locus, there were three variants in or near *CYP2A6* associated with CPD

independently of the established variant (rs7937) and each other: rs8102683 (conditional  $P = 4.53 \times 10^{-16}$ ), rs28399442 (conditional  $P = 2.63 \times 10^{-12}$ ) and rs3865453 (conditional  $P = 4.96 \times 10^{-10}$ ) and rs28399442 was a low-frequency variant. The same SNVs also showed evidence of independent effects with pack-years, albeit with larger  $P$ -values ( $P < 5 \times 10^{-6}$ ; Supp. Table 5). At the *TEX41/PABPC1P2* locus, rs11694518 (conditional  $P = 3.43 \times 10^{-7}$ ) was associated with smoking initiation independently of the established variant (rs10427255). At 15q25, rs938682 ( $P = 7.78 \times 10^{-21}$ ) was associated with CPD independently of the established variant (rs1051730) and (in agreement with a previous report [49]) is an eQTL for *CHRNA5* in brain putamen basal ganglia tissues in GTEx.

### Gene-based association studies

Gene-based collapsing tests using  $MAF < 0.01$  variants, did not identify any associated genes at the pre-specified  $P < 8 \times 10^{-7}$  threshold. Of the top four gene associations, three were novel (*CHRNA2*, *MMP17*, and *CRCP*) and one was known (*CHRNA5*), and had  $P < 7 \times 10^{-4}$ , with CPD and/or pack-years (Supp. Table 6). Analyses conditional on the variant with the smallest  $P$ -value in the gene, revealed the associations at *CHRNA2*, *MMP17* and *CRCP* were due to more than one rare variant (conditional  $P < 0.05$ ; Supp. Table 6). In contrast, the *CHRNA5* gene association was attributable to a single variant (rs2229961).

### Mendelian randomization analyses

We conducted MR analyses to elucidate the potential causal impact of SI and CPD on BMI, schizophrenia and EA using the MR-Egger, median weighted and inverse variance weighted methods. We found a causal association between SI and EA using both the median weighted and inverse variance weighted methods ( $P < 0.0001$ ; Supp. Table 7) but not with MR-Egger ( $P = 0.2$ ). There was an association of SI with BMI using MR-Egger only ( $P = 0.01$ ; Supp. Table 7), but there was evidence of horizontal pleiotropy ( $P = 0.001$ ) and no support from the other methods. Similarly, increased CPD was only associated with reduced BMI using the weighted median approach ( $P = 0.009$ ) and not the other methods ( $P > 0.017$ ). We also tested if schizophrenia, EA or BMI causally influence CPD or SI using SNVs associated with schizophrenia, EA and BMI, respectively, as instrumental variables. No evidence of such reverse causation was found (Supp. Table 7). These results were consistent with previous analyses [50]. There was no evidence of a causal effect of SI on schizophrenia, or CPD on educational attainment (Supp. Table 7).

### Functional characterization of novel loci

Using proxies with  $r^2 \geq 0.8$  in 1000 Genomes Europeans, we queried the GWAS catalogue [38] ( $P \leq 5 \times 10^{-8}$ ) for pleiotropic effects of our novel sentinel SNVs. Two, rs11539157 and rs3001723 were previously associated with schizophrenia [36], suggesting shared biological pathways between schizophrenia and smoking behaviours (Table 2). This fits with the known association of smoking with schizophrenia [51]. Two, rs1514175 and rs2947411 have previously been associated with BMI [52], and extreme obesity [53].

eQTL lookups in GTEx V7 (13 Brain tissues with  $\geq 80$  samples) [39], Brain xQTL [40] and BRAINEAC [41] databases revealed that the A allele at rs462779, which decreases risk of smoking initiation, also decreased expression of *REV3L* in cerebellum in GTEx (A allele  $P = 4.8 \times 10^{-8}$ ;  $\beta = -0.40$ ) and was in strong LD with the top eQTL for *REV3L* in cerebellum ( $r^2 = 0.86$  with rs9487668 in 1000 Genomes Europeans). The smoking initiation-associated SNV, rs12780116, was an eQTL for *BORCS7* in four brain tissues, and *NT5C2* in the cerebellar hemisphere (A allele  $P = 4.5 \times 10^{-7}$ ;  $\beta = -0.32$ ) and the cerebellum ( $P = 5.6 \times 10^{-6}$ ;  $\beta = -0.415$ ; in strong LD with the top eQTL,  $r^2 = 0.97$  with rs11191546). The G allele of a second variant in the region, rs7096169 (intronic to *BORCS7* and only in weak LD with rs12780116,  $r^2 = 0.18$  in 1000G Europeans) increases smoking initiation and reduces expression of *BORCS7* and *AS3MT* in eight brain tissues (including dorsolateral prefrontal cortex in the Brain xQTL and was the top *BORCS7* eSNP in GTEx in the Cerebellar Hemisphere, Cerebellum, and Spinal cord cervical-C1). The same variant also reduced expression of *ARL3* in cerebellum in GTEx (Table 2).

Biological pathway enrichment analyses carried out in ConsensusPathDB [44] using the genes implicated by the eQTL databases (Table 2) and/or a coding SNVs (i.e., *PJA1*, *GPR101*) showed that the (i) pyrimidine metabolism and (ii) activation of nicotinic acetylcholine receptors pathways are enriched for these smoking behaviour associated genes (false discovery rate  $< 0.01$ ;  $P < 0.0001$ ).

### Discussion

Smoking is the most important preventable lifestyle risk factor for many diseases, including cancers [3, 54], heart disease [4, 55] and many respiratory diseases such as COPD [1, 2]. Not initiating is the best way to prevent smoking-related diseases and genetics can play a considerable part in smoking behaviours including initiation. We have performed the largest exome-wide genetic association study of smoking behaviour-related traits to date involving up to

622,409 individuals, and identified and replicated five associations, including two on the X-chromosome (Table 1). We identified a further 35 novel associations in a meta-analysis of discovery and replication cohorts (Table 2). We validated 14 previously reported SNV-smoking trait associations (Supp. Table 3) and identified secondary independent associations at three loci, including three in the 19q13 region (rs8102683, rs28399442, and rs3865453; Table 3).

Gene-based tests improve power by aggregating effects of rare variants. While no genes reached our Bonferroni-adjusted  $P$ -value threshold, we identified three candidate genes with multiple rare variant associations for future replication: calcitonin gene-related peptide-receptor component (*CRCP*) with CPD and *CHRNA2* and *MMP17* with pack-years (Supp. Table 6; also see ‘Genes of Interest’ section in Supp. Material). *CRCP*’s protein product is expressed in brain tissues amongst others and functions as part of a receptor complex for a neuropeptide that increases intracellular cyclic adenosine monophosphate levels [56]. *MMP17* encodes a matrix metalloproteinase that is also expressed in the brain and is a member of the peptidase M10 family, and proteins in this family are involved in the breakdown of extracellular matrix in normal physiological processes [57]. Given, we were not able conclusively to identify rare variant associations, even larger studies, are required to identify rare variants associated with smoking behaviours. In addition, phenotypes such as cotinine levels [58] and nicotine metabolism speed [59] could be interrogated using methods such as MTAG [60] to improve power.

As recommended by UK Biobank, we analysed UK Biobank samples by adjusting for genotyping array because a subset of (extreme smokers in) UK Biobank were genotyped on a different array (UK BiLEVE). However, this adjustment could potentially introduce collider bias in analyses of smoking traits. Given that the UK BiLEVE study is relatively small compared to the full study, and the genetic effect sizes for smoking-associated variants are small, we expect the influence of collider bias to be small [61]. Nevertheless, we performed sensitivity analyses to assess the impact of collider bias. Firstly, we performed a meta-analysis excluding the UK BiLEVE samples, and secondly, we re-analysed UK Biobank without adjusting for genotype array. As expected, the estimated genetic effects from these additional analyses were very similar to our reported results suggesting collider bias is not a concern (Suppl. Table 8).

Follow-up of the replicated SNVs in the literature and eQTL databases implicated some potentially interesting genes: *NT5C2* is known to hydrolyse purine nucleotides and be involved in maintaining cellular nucleotide balance, and was previously associated with schizophrenia [62]. *REV3L*, encodes the catalytic subunit of DNA polymerase  $\zeta$  (zeta) which is involved in translesion DNA synthesis. Previously, polymorphisms in a microRNA target site of *REV3L* were

shown to be associated with lung cancer susceptibility [63]. We showed that decreased expression of *REV3L* may also lower the probability of smoking initiation. The SNV, rs11776293, intronic in *EPHX2*, was associated with reduced SI in the combined meta-analysis, and is in LD with rs56372821 ( $r^2 = 0.83$ ), which is associated with reduced cannabis use disorder [64]. rs216195 (in *SMG6*) was genome-wide significant in the discovery and the combined meta-analysis. *SMG6* is a plausible candidate gene as it was previously shown to be less methylated in current smokers compared to never smokers [65]. The combined meta-analysis also identified a rare missense variant in *CCDC141*, rs150493199 (MAF < 0.01; Table 2). Coding variants in *CCDC141* were previously associated with heart rate [66] and blood pressure [67, 68].

Smoking behaviours represent a complex phenotype that are linked to an array of socio-cultural and familial, as well as genetic determinants. Kong et al., recently reported that ‘genetic-nurture’ i.e., effects of non-transmitted parental alleles, affect educational attainment [69]. They also show that there is an effect of educational attainment and genetic nurture on smoking behaviour. Four of our sentinel SNVs (or a strong proxy;  $r^2 > 0.8$ ) were associated with years of educational attainment [37] (rs2292239, rs3001723 ( $P < 5 \times 10^{-8}$ ), rs9320995 ( $P = 8.90 \times 10^{-7}$ ), and rs13022438 ( $P = 3.79 \times 10^{-6}$ ), in agreement with this paradigm and our MR analyses indicated that initiating smoking reduced years in education. Future family studies will be required to disentangle how much of the variance explained in the current analysis is due to direct versus genetic nurturing effects.

Our study primarily focused on European ancestry, but we also included two non-European studies but these non-European studies lacked statistical power on their own to identify ancestry-specific effects. Therefore, we did not perform ancestry-specific meta-analyses. Nevertheless, our results offered cross ancestry replication. One of the associations identified in the conditional analyses, rs8102683 (near *CYP2A6*), confirmed an association with CPD that was previously identified by Kumasaka et al. in a Japanese population [70] but this is the first time it was associated in Europeans (rs8102683 is also correlated with rs56113850 ( $r^2 = 0.43$ ), a SNV identified previously by Loukola et al. [59] in a genetic association study of nicotine metabolite ratio in Europeans). As more non-European studies become available, it would be of great interest to perform non-European ancestry studies, in order to fine-map causal variants for smoking-related traits.

CPD and pack-years are two correlated measures of smoking. In the ~40,000 individuals from UK Biobank with CPD and pack-years calculated, correlation between CPD and pack-years was 0.640. Interestingly, while pack-years was inversely correlated with smoking cessation (−0.18) i.e., the more years a smoker has been smoking the less



likely they were to cease, CPD was positively correlated with smoking cessation (0.13) i.e., heavier smokers were more likely to stop smoking. In contrast, the *DBH* SNV, rs3025343, (first identified via its association with increased smoking cessation [6]) was associated with increased pack-years ( $P = 1.29 \times 10^{-14}$ ) and increased CPD ( $P = 2.93 \times 10^{-9}$ ) in our study. The association at *DBH* also represents the first time that a SNV has a smaller *P*-value for pack-years ( $n = 131,892$ ) compared to CPD ( $n = 128,746$ ). These findings may help elucidate the genetic basis of these correlated addiction phenotypes.

We performed the largest exome-wide genetic association study of smoking behaviour-related traits to date and nearly doubled the number of replicated associations to 24 (including conditional analyses) including associations on the X-chromosome for the first time, which merit further study. We also identified a further 35 novel smoking trait associated SNVs in the combined meta-analysis. The novel loci identified in this study will substantially expand our knowledge of the smoking addiction-related traits, facilitate understanding the genetic aetiology of smoking behaviour and may lead to the identification of drug targets of potential relevance to prevent individuals from initiating smoking and/or aid smokers to stop smoking.

**Acknowledgments** Author acknowledgements are included in the Supplementary material.

## Compliance with ethical standards

**Conflict of interest** Paul W. Franks has been a paid consultant for Eli Lilly and Sanofi Aventis and has received research support from several pharmaceutical companies as part of European Union Innovative Medicines Initiative (IMI) projects. Neil Poulter has received financial support from several pharmaceutical companies that manufacture either blood pressure lowering or lipid lowering agents or both and consultancy fees. Peter Sever has received research awards from Pfizer. Mark J. Caulfield is Chief Scientist for Genomics England, a UK government company. Adam Butterworth reports grants from Merck and Pfizer and grants outside of this work from AstraZeneca, Biogen, and Novartis, and personal fees from Novartis. John Danesh reports grants, personal fees and non-financial support from Merck Sharp & Dohme (MSD), grants, personal fees and non-financial support from Novartis, grants from British Heart Foundation, grants from European Research Council, grants from NIHR, grants from NHS Blood and Transplant, grants from Pfizer, grants from UK MRC, grants from Wellcome Trust, grants from AstraZeneca, outside the submitted work. Laura J. Bierut is listed as an inventor on Issued U.S. Patent 8,080,371, "Markers for Addiction" covering the use of certain SNPs in determining the diagnosis, prognosis, and treatment of addiction.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless

indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

1. Wain LV, Shrine N, Miller S, Jackson VE, Ntalla I, Soler Artigas M, et al. Novel insights into the genetics of smoking behaviour, lung function, and chronic obstructive pulmonary disease (UK BiLEVE): a genetic association study in UK Biobank. *Lancet Respir Med*. 2015;3:769–81.
2. Wain LV, Shrine N, Artigas MS, Erzurumluoglu AM, Noyvert B, Bossini-Castillo L, et al. Genome-wide association analyses for lung function and chronic obstructive pulmonary disease identify new loci and potential druggable targets. *Nat Genet*. 2017;49:416–25.
3. McKay JD, Hung RJ, Han Y, Zong X, Carreras-Torres R, Christiani DC, et al. Large-scale association analysis identifies new lung cancer susceptibility loci and heterogeneity in genetic susceptibility across histological subtypes. *Nat Genet*. 2017;49:1126–32.
4. O'Donnell CJ, Nabel EG. Genomics of Cardiovascular Disease. *New Engl J Med*. 2011;365:2098–109.
5. Reitsma MB, Fullman N, Ng M, Salama JS, Abajobir A, Abate KH, et al. Smoking prevalence and attributable disease burden in 195 countries and territories, 1990–2015: a systematic analysis from the Global Burden of Disease Study 2015. *The Lancet* 2017;389:1885–906.
6. Tobacco and Genetics Consortium. Genome-wide meta-analyses identify multiple loci associated with smoking behavior. *Nat Genet*. 2010;42:441–7.
7. Hancock DB, Reginsson GW, Gaddis NC, Chen X, Saccone NL, Lutz SM, et al. Genome-wide meta-analysis reveals common splice site acceptor variant in *CHRNA4* associated with nicotine dependence. *Transl Psychiatry*. 2015;5:e651.
8. Siedlinski M, Cho MH, Bakke P, Gulsvik A, Lomas DA, Anderson W, et al. Genome-wide association study of smoking behaviours in patients with COPD. *Thorax*. 2011;66:894–902.
9. Thorgeirsson TE, Gudbjartsson DF, Surakka I, Vink JM, Amin N, Geller F, et al. Sequence variants at *CHRNA3-CHRNA6* and *CYP2A6* affect smoking behavior. *Nat Genet*. 2010;42:448–53.
10. Timofeeva MN, McKay JD, Smith GD, Johansson M, Byrnes GB, Chabrier A, et al. Genetic polymorphisms in 15q25 and 19q13 loci, cotinine levels, and risk of lung cancer in EPIC. *Cancer Epidemiol Biomark Prev*. 2011;20:2250–61.
11. Bloom AJ, Baker TB, Chen L-S, Breslau N, Hatsukami D, Bierut LJ, et al. Variants in two adjacent genes, *EGLN2* and *CYP2A6*, influence smoking behavior related to disease risk via different mechanisms. *Hum Mol Genet*. 2014;23:555–61.
12. Thakur GA, Sengupta SM, Grizenko N, Choudhry Z, Joober R. Family-based association study of ADHD and genes increasing the risk for smoking behaviours. *Arch Dis Child*. 2012;97:1027.
13. Munafò MR, Flint J. The genetic architecture of psychophysiological phenotypes. *Psychophysiology*. 2014;51:1331–2.
14. Keskitalo K, Broms U, Heliövaara M, Ripatti S, Surakka I, Perola M, et al. Association of serum cotinine level with a cluster of three nicotinic acetylcholine receptor genes (*CHRNA3/CHRNA5/CHRNA4*) on chromosome 15. *Hum Mol Genet*. 2009;18:4007–12.
15. Vink JM, Willemsen G, Boomsma DI. Heritability of smoking initiation and nicotine dependence. *Behav Genet*. 2005;35:397–406.
16. Carmelli D, Swan GE, Robinette D, Fabsitz R. Genetic Influence on Smoking — A Study of Male Twins. *New Engl J Med*. 1992;327:829–33.

17. Kaprio J, Koskenvuo M, Sarna S. Cigarette smoking, use of alcohol, and leisure-time physical activity among same-sexed adult male twins. *Prog Clin Biol Res.* 1981;69(Pt C):37–46.
18. Liu DJ, Brazel DM, Turcot V, Zhan X, Gong J, Barnes DR, et al. Exome chip meta-analysis elucidates the genetic architecture of rare coding variants in smoking and drinking behavior. *bioRxiv* 2017.
19. Centers for Disease Control and Prevention (CDC). Cigarette smoking among adults—United States, 2007. *MMWR Morb Mortal Wkly Rep.* 2008;57:1221–6.
20. Staley JR, Jones E, Kaptoge S, Butterworth AS, Sweeting MJ, Wood AM, et al. A comparison of Cox and logistic regression for use in genome-wide association studies of cohort and case-cohort design. *Eur J Human Genet : EJHG.* 2017;25:854–62.
21. Feng S, Liu D, Zhan X, Wing MK, Abecasis GR. RAREMETAL: fast and powerful meta-analysis for rare variants. *Bioinformatics.* 2014;30:2828–9.
22. Zhan X, Hu Y, Li B, Abecasis GR, Liu DJ. RVTESTS: an efficient and comprehensive tool for rare variant association analysis using sequence data. *Bioinformatics.* 2016;32:1423–6.
23. Willer CJ, Li Y, Abecasis GR. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics.* 2010;26:2190–1.
24. Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, et al. The UK Biobank resource with deep phenotyping and genomic data. *Nature.* 2018;562:203–9.
25. Jiang B, Chen S, Jiang Y, Liu M, Iacono WG, Hewitt JK, et al. Proper Conditional Analysis in the Presence of Missing Data Identified Novel Independently Associated Low Frequency Variants in Nicotine Dependence Genes. *bioRxiv* 2017.
26. Madsen BE, Browning SR. A groupwise association test for rare mutations using a weighted sum statistic. *PLoS Genet.* 2009;5:e1000384.
27. Morris AP, Zeggini E. An evaluation of statistical approaches to rare variant analysis in genetic association studies. *Genet Epidemiol.* 2010;34:188–93.
28. Wu MC. Rare variant association testing for sequencing data using the sequence kernel association test (SKAT). *Am J Hum Genet.* 2011;89:82–93.
29. Zhan X, Liu DJ. SEQMINER: An R-Package to Facilitate the Functional Interpretation of Sequence-Based Associations. *Genet Epidemiol.* 2015;39:619–23.
30. Hemani G, Zheng J, Elsworth B, Wade KH, Haberland V, Baird D, et al. The MR-Base platform supports systematic causal inference across the human phenome. *eLife.* 2018;7:e34408.
31. Pierce BL, Burgess S. Efficient design for Mendelian randomization studies: subsample and 2-sample instrumental variable estimators. *Am J Epidemiol.* 2013;178:1177–84.
32. Bowden J, Davey Smith G, Burgess S. Mendelian randomization with invalid instruments: effect estimation and bias detection through Egger regression. *Int J Epidemiol.* 2015;44:512–25.
33. Rees JMB, Wood AM, Burgess S. Extending the MR-Egger method for multivariable Mendelian randomization to correct for both measured and unmeasured pleiotropy. *Stat Med.* 2017;36:4705–18.
34. Bowden J, Davey Smith G, Haycock PC, Burgess S. Consistent estimation in Mendelian randomization with some invalid instruments using a weighted median estimator. *Genet Epidemiol.* 2016;40:304–14.
35. Locke AE, Kahali B, Berndt SI, Justice AE, Pers TH, Day FR, et al. Genetic studies of body mass index yield new insights for obesity biology. *Nature.* 2015;518:197–206.
36. Schizophrenia Working Group of the Psychiatric Genomics Consortium, Ripke S, Neale BM, Corvin A, Walters JTR, Farh K-H, et al. Biological insights from 108 schizophrenia-associated genetic loci. *Nature.* 2014;511:421.
37. Okbay A, Beauchamp JP, Fontana MA, Lee JJ, Pers TH, Rietveld CA, et al. Genome-wide association study identifies 74 loci associated with educational attainment. *Nature.* 2016;533:539–42.
38. MacArthur J, Bowler E, Cerezo M, Gil L, Hall P, Hastings E, et al. The new NHGRI-EBI Catalog of published genome-wide association studies (GWAS Catalog). *Nucleic Acids Res.* 2017;45(D1):D896–901.
39. Battle A, Brown CD, Engelhardt BE, Montgomery SB. Genetic effects on gene expression across human tissues. *Nature.* 2017;550:204–13.
40. Ng B, White CC, Klein H-U, Sieberts SK, McCabe C, Patrick E, et al. An xQTL map integrates the genetic architecture of the human brain's transcriptome and epigenome. *Nat Neurosci* 2017; 20:1418–26; advance online publication.
41. Trabzuni D, Ryten M, Walker R, Smith C, Imran S, Ramasamy A, et al. Quality control parameters on a large dataset of regionally dissected human control brains for whole genome expression studies. *J Neurochem.* 2011;119:275–82.
42. Ongen H, Buil A, Brown AA, Dermitzakis ET, Delaneau O. Fast and efficient QTL mapper for thousands of molecular phenotypes. *Bioinformatics.* 2016;32:1479–85.
43. Storey JD, Tibshirani R. Statistical significance for genomewide studies. *Proc Natl Acad Sci USA.* 2003;100:9440–5.
44. Kamburov A, Wierling C, Lehrach H, Herwig R. ConsensusPathDB—a database for integrating human functional interaction networks. *Nucleic Acids Res.* 2009;37(suppl\_1):D623–8.
45. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, et al. A method and server for predicting damaging missense mutations. *Nat Meth.* 2010;7:248–9.
46. Shihab HA, Gough J, Cooper DN, Stenson PD, Barker GL, Edwards KJ, et al. Predicting the functional, molecular, and phenotypic consequences of amino acid substitutions using hidden Markov models. *Hum Mutat.* 2013;34:57–65.
47. Bulik-Sullivan BK, Loh PR, Finucane HK, Ripke S, Yang J, Patterson N, et al. LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat Genet.* 2015;47:291–5.
48. Zheng J, Erzurumluoglu AM, Elsworth BL, Kemp JP, Howe L, Haycock PC, et al. LD Hub: a centralized database and web interface to perform LD score regression that maximizes the potential of summary level GWAS data for SNP heritability and genetic correlation analysis. *Bioinformatics.* 2017;33:272–9.
49. Wang JC, Cruchaga C, Saccone NL, Bertelsen S, Liu P, Budde JP, et al. Risk for nicotine dependence and lung cancer is conferred by mRNA expression levels and amino acid change in CHRNA5. *Hum Mol Genet.* 2009;18:3125–35.
50. Gage SH, Jones HJ, Taylor AE, Burgess S, Zammit S, Munafò MR. Investigating causality in associations between smoking initiation and schizophrenia using Mendelian randomization. *Sci Rep.* 2017;7:40653.
51. Kelly C, McCreddie R. Cigarette smoking and schizophrenia. *Adv Psychiatr Treat.* 2000;6:327–31.
52. Speliotes EK, Willer CJ, Berndt SI, Monda KL, Thorleifsson G, Jackson AU, et al. Association analyses of 249,796 individuals reveal 18 new loci associated with body mass index. *Nat Genet.* 2010;42:937–48.
53. Wheeler E, Huang N, Bochukova EG, Keogh JM, Lindsay S, Garg S, et al. Genome-wide SNP and CNV analysis identifies common and low-frequency variants associated with severe early-onset obesity. *Nat Genet.* 2013;45:513–7.
54. Hecht SS. Tobacco Smoke Carcinogens and Lung Cancer. *J Natl Cancer Inst.* 1999;91:1194–210.
55. Ockene IS, Miller NH. Cigarette Smoking, Cardiovascular Disease, and Stroke. A Statement Healthc Prof Am Heart Assoc. 1997;96:3243–7.



56. Uhlen M, Fagerberg L, Hallstrom BM, Lindskog C, Oksvold P, Mardinoglu A, et al. Proteomics. Tissue-based map of the human proteome. *Science*. 2015;347:1260419.
57. O'Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, et al. Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res*. 2016;44(D1):D733–45.
58. Ware JJ, Chen X, Vink J, Loukola A, Minica C, Pool R, et al. Genome-Wide Meta-Analysis of Cotinine Levels in Cigarette Smokers Identifies Locus at 4q13.2. *Sci Rep*. 2016;6:20092.
59. Loukola A, Buchwald J, Gupta R, Palviainen T, Hallfors J, Tikkanen E, et al. A Genome-Wide Association Study of a Biomarker of Nicotine Metabolism. *PLoS Genet*. 2015;11:e1005498.
60. Turley P, Walters RK, Maghizian O, Okbay A, Lee JJ, Fontana MA, et al. Multi-trait analysis of genome-wide association summary statistics using MTAG. *Nat Genet*. 2018;50:229–37.
61. Munafo MR, Tilling K, Taylor AE, Evans DM, Davey Smith G. Collider scope: when selection bias can substantially influence observed associations. *Int J Epidemiol*. 2018;47:226–35.
62. Aberg KA, Liu Y, Bukszar J, et al. A comprehensive family-based replication study of schizophrenia genes. *JAMA Psychiatry*. 2013;70:573–81.
63. Zhang S, Chen H, Zhao X, Cao J, Tong J, Lu J, et al. REV3L 3'UTR 460 T>C polymorphism in microRNA target sites contributes to lung cancer susceptibility. *Oncogene*. 2013;32:242–50.
64. Demontis D, Rajagopal VM, Als TD, Grove J, Pallesen J, Hjerthoj C, et al. Genome-wide association study implicates *CHRNA2* in cannabis use disorder. *bioRxiv* 2018.
65. Steenaard RV, Ligthart S, Stolk L, Peters MJ, van Meurs JB, Uitterlinden AG, et al. Tobacco smoking is associated with methylation of genes related to coronary artery disease. *Clin Epigenetics*. 2015;7:54.
66. van den Berg ME, Warren HR, Cabrera CP, Verweij N, Mifsud B, Haessler J, et al. Discovery of novel heart rate-associated loci using the Exome Chip. *Hum Mol Genet*. 2017;26:2346–63.
67. Warren HR, Evangelou E, Cabrera CP, Gao H, Ren M, Mifsud B, et al. Genome-wide association analysis identifies novel blood pressure loci and offers biological insights into cardiovascular risk. *Nat Genet*. 2017;49:403–15.
68. Hoffmann TJ, Ehret GB, Nandakumar P, Ranatunga D, Schaefer C, Kwok PY, et al. Genome-wide association analyses using electronic health records identify new loci influencing blood pressure variation. *Nat Genet*. 2017;49:54–64.
69. Kong A, Thorleifsson G, Frigge ML, Vilhjalmsdottir BJ, Young AI, Thorgeirsson TE, et al. The nature of nurture: effects of parental genotypes. *Science*. 2018;359:424–8.
70. Kumasaka N, Aoki M, Okada Y, Takahashi A, Ozaki K, Mushi-iroda T, et al. Haplotypes with copy number and single nucleotide polymorphisms in CYP2A6 locus are associated with smoking quantity in a Japanese population. *PLoS ONE* 2012;7:e44507.

## Author List

A. Mesut Erzurumluoglu<sup>1</sup> · Mengzhen Liu<sup>2</sup> · Victoria E. Jackson<sup>1,3,4</sup> · Daniel R. Barnes<sup>5</sup> · Gargi Datta<sup>2,6</sup> · Carl A. Melbourne<sup>1</sup> · Robin Young<sup>5</sup> · Chiara Batini<sup>1</sup> · Praveen Surendran<sup>5</sup> · Tao Jiang<sup>5</sup> · Sheikh Daud Adnan<sup>7</sup> · Saima Afaq<sup>8</sup> · Arpana Agrawal<sup>9</sup> · Elisabeth Altmaier<sup>10</sup> · Antonis C. Antoniou<sup>11</sup> · Folkert W. Asselbergs<sup>12,13,14,15</sup> · Clemens Baumbach<sup>10</sup> · Laura Bierut<sup>16</sup> · Sarah Bertelsen<sup>17</sup> · Michael Boehnke<sup>18</sup> · Michiel L. Bots<sup>19,20</sup> · David M Brazel<sup>6,21</sup> · John C. Chambers<sup>8,22,23,24</sup> · Jenny Chang-Claude<sup>25,26</sup> · Chu Chen<sup>27,28</sup> · Janie Corley<sup>29,30</sup> · Yi-Ling Chou<sup>9</sup> · Sean P. David<sup>31</sup> · Rudolf A. de Boer<sup>32</sup> · Christiaan A. de Leeuw<sup>33</sup> · Joe G. Dennis<sup>11</sup> · Anna F. Dominiczak<sup>34</sup> · Alison M. Dunning<sup>35</sup> · Douglas F. Easton<sup>11,35</sup> · Charles Eaton<sup>28</sup> · Paul Elliott<sup>36,37,38,39</sup> · Evangelos Evangelou<sup>8,40</sup> · Jessica D. Faul<sup>80</sup> · Tatiana Foroud<sup>41</sup> · Alison Goate<sup>42</sup> · Jian Gong<sup>43</sup> · Hans J. Grabe<sup>44</sup> · Jeff Haessler<sup>43</sup> · Christopher Haiman<sup>45</sup> · Göran Hallmans<sup>46</sup> · Anke R. Hammerschlag<sup>33</sup> · Sarah E. Harris<sup>29,47</sup> · Andrew Hattersley<sup>48</sup> · Andrew Heath<sup>9</sup> · Chris Hsu<sup>49</sup> · William G. Iacono<sup>2</sup> · Stavroula Kanoni<sup>50,51</sup> · Manav Kapoor<sup>17</sup> · Jaakko Kaprio<sup>52,53</sup> · Sharon L. Kardia<sup>54</sup> · Fredrik Karpe<sup>55,56</sup> · Jukka Kontto<sup>57</sup> · Jaspal S. Kooner<sup>23,24,37,58</sup> · Charles Kooperberg<sup>43,59</sup> · Kari Kuulasmaa<sup>57</sup> · Markku Laakso<sup>60</sup> · Dongbing Lai<sup>41</sup> · Claudia Langenberg<sup>61</sup> · Nhung Le<sup>62</sup> · Guillaume Lettre<sup>63,64</sup> · Anu Loukola<sup>52,53</sup> · Jian'an Luan<sup>61</sup> · Pamela A. F. Madden<sup>9</sup> · Massimo Mangino<sup>65,129</sup> · Riccardo E. Marioni<sup>29,47</sup> · Eirini Marouli<sup>50,51</sup> · Jonathan Marten<sup>66</sup> · Nicholas G. Martin<sup>67</sup> · Matt McGue<sup>2</sup> · Kyriaki Michailidou<sup>11,68</sup> · Evelin Mihailov<sup>69</sup> · Alireza Moayyeri<sup>70</sup> · Marie Moitry<sup>71</sup> · Martina Müller-Nurasyid<sup>72,73,74</sup> · Aliya Naheed<sup>75</sup> · Matthias Nauck<sup>76,77</sup> · Matthew J. Neville<sup>55,56</sup> · Sune Fallgaard Nielsen<sup>78</sup> · Kari North<sup>79</sup> · Markus Perola<sup>52,57</sup> · Paul D. P. Pharoah<sup>11,35</sup> · Giorgio Pistis<sup>81</sup> · Tinca J. Polderman<sup>33</sup> · Danielle Posthuma<sup>33,82</sup> · Neil Poulter<sup>83</sup> · Beenish Qaiser<sup>52,53</sup> · Asif Rasheed<sup>84</sup> · Alex Reiner<sup>28,43</sup> · Frida Renström<sup>85,86</sup> · John Rice<sup>87</sup> · Rebecca Rohde<sup>88</sup> · Olov Rolandsson<sup>89</sup> · Nilesh J. Samani<sup>90</sup> · Maria Samuel<sup>84</sup> · David Schlessinger<sup>91</sup> · Steven H Scholte<sup>92</sup> · Robert A. Scott<sup>61</sup> · Peter Sever<sup>58,83</sup> · Yaming Shao<sup>88</sup> · Nick Shrine<sup>1</sup> · Jennifer A. Smith<sup>54</sup> · John M. Starr<sup>29,93</sup> · Kathleen Stirrups<sup>50,94</sup> · Danielle Stram<sup>95</sup> · Heather M. Stringham<sup>18</sup> · Ioanna Tachmazidou<sup>96</sup> · Jean-Claude Tardif<sup>63,64</sup> · Deborah J. Thompson<sup>11</sup> · Hilary A. Tindle<sup>97</sup> · Vinicius Tragante<sup>98</sup> · Stella Trompet<sup>99,100</sup> · Valerie Turcot<sup>63</sup> · Jessica Tyrrell<sup>48</sup> · Ilonca Vaartjes<sup>19,20</sup> · Andries R van der Leij<sup>92</sup> · Peter van der Meer<sup>32</sup> · Tibor V. Varga<sup>85</sup> · Niek Verweij<sup>32,101</sup> · Henry Völzke<sup>77,102</sup> · Nicholas J. Wareham<sup>61</sup> · Helen R. Warren<sup>103,104</sup> · David R. Weir<sup>80</sup> · Stefan Weiss<sup>77,105</sup> · Leah Wetherill<sup>41</sup> · Hanieh Yaghootkar<sup>48</sup> · Ersin Yavas<sup>106,107</sup> · Yu Jiang<sup>108</sup> · Fang Chen<sup>108</sup> · Xiaowei Zhan<sup>109</sup> · Weihua Zhang<sup>8,110</sup> · Wei Zhao<sup>111</sup> · Wei Zhao<sup>54</sup> · Kaixin Zhou<sup>112</sup> · Philippe Amouyel<sup>113</sup> · Stefan Blankenberg<sup>114,115</sup> · Mark J. Caulfield<sup>103,104</sup> · Rajiv Chowdhury<sup>5</sup> · Francesco Cucca<sup>81</sup> · Ian J. Deary<sup>29,30</sup> · Panos Deloukas<sup>96,116,117</sup> · Emanuele Di Angelantonio<sup>5,118</sup> · Marco Ferrario<sup>119</sup>

Jean Ferrières<sup>120</sup> · Paul W. Franks<sup>85,121</sup> · Tim M. Frayling<sup>48</sup> · Philippe Frossard<sup>84</sup> · Ian P. Hall<sup>122</sup> · Caroline Hayward<sup>66</sup> · Jan-Håkan Jansson<sup>123</sup> · J. Wouter Jukema<sup>124,125</sup> · Frank Kee<sup>126</sup> · Satu Männistö<sup>57</sup> · Andres Metspalu<sup>69</sup> · Patricia B. Munroe<sup>103,104</sup> · Børge Grønne Nordestgaard<sup>78</sup> · Colin N. A. Palmer<sup>127</sup> · Veikko Salomaa<sup>57</sup> · Naveed Sattar<sup>128</sup> · Timothy Spector<sup>129</sup> · David Peter Strachan<sup>130</sup> · Understanding Society Scientific Group, EPIC-CVD, GSCAN, Consortium for Genetics of Smoking Behaviour, CHD Exome+ consortium · Pim van der Harst<sup>32,131</sup> · Eleftheria Zeggini<sup>96</sup> · Danish Saleheen<sup>5,132,133</sup> · Adam S. Butterworth<sup>5,118</sup> · Louise V. Wain<sup>1,134</sup> · Goncalo R. Abecasis<sup>18</sup> · John Danesh<sup>5,96,118</sup> · Martin D. Tobin<sup>1,134</sup> · Scott Vrieze<sup>2</sup> · Dajiang J. Liu<sup>108</sup> · Joanna M. M. Howson<sup>id 5</sup>

<sup>1</sup> Department of Health Sciences, University of Leicester, Leicester, UK

<sup>2</sup> Department of Psychology, University of Minnesota, Minneapolis, MN, USA

<sup>3</sup> Population Health and Immunity Division, The Walter and Eliza Hall Institute of Medical Research, 1G Royal Pde, 3052 Parkville, Australia

<sup>4</sup> Department of Medical Biology, University of Melbourne, Melbourne 3010 Parkville, Australia

<sup>5</sup> MRC/BHF Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of Cambridge, Cambridge, CB1 8RN, UK

<sup>6</sup> Institute for Behavioral Genetics, University of Colorado Boulder, Boulder, CO, USA

<sup>7</sup> National Institute of Cardiovascular Diseases, Sher-e-Bangla Nagar, Dhaka, Bangladesh

<sup>8</sup> Department of Epidemiology and Biostatistics, Imperial College London, London W2 1PG, UK

<sup>9</sup> Department of Psychiatry, Washington University, St. Louis, MO, USA

<sup>10</sup> Research Unit of Molecular Epidemiology, Helmholtz Zentrum München-German Research Center for Environmental Health, Neuherberg, Germany

<sup>11</sup> Centre for Cancer Genetic Epidemiology, Department of Public Health and Primary Care, University of Cambridge, Cambridge CB1 8RN, UK

<sup>12</sup> Department of Cardiology, Division Heart & Lungs, University Medical Center Utrecht, University of Utrecht, Utrecht, The Netherlands

<sup>13</sup> Durrer Center for Cardiovascular Research, Netherlands Heart Institute, Utrecht, The Netherlands

<sup>14</sup> Institute of Cardiovascular Science, Faculty of Population Health Sciences, University College London, London, UK

<sup>15</sup> Farr Institute of Health Informatics Research and Institute of Health Informatics, University College London, London, UK

<sup>16</sup> Department of Psychiatry, Washington University School of Medicine, St. Louis, MO, USA

<sup>17</sup> Department of Neuroscience, Icahn School of Medicine at Mount Sinai, New York, NY, USA

<sup>18</sup> Department of Biostatistics and Center for Statistical Genetics, University of Michigan, Ann Arbor, MI, USA

<sup>19</sup> Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, 3508GA Utrecht, The Netherlands

<sup>20</sup> Center for Circulatory Health, University Medical Center Utrecht, 3508GA Utrecht, The Netherlands

<sup>21</sup> Department of Molecular, Cellular, and Developmental Biology, University of Colorado Boulder, Boulder, CO, USA

<sup>22</sup> Lee Kong Chian School of Medicine, Nanyang Technological University, Singapore 308232, Singapore

<sup>23</sup> Department of Cardiology, Ealing Hospital, Middlesex UB1 3HW, UK

<sup>24</sup> Imperial College Healthcare NHS Trust, London W12 0HS, UK

<sup>25</sup> Division of Cancer Epidemiology, German Cancer Research Centre (DKFZ), Heidelberg, Germany

<sup>26</sup> Cancer Epidemiology Group, University Medical Centre Hamburg-Eppendorf, University Cancer Centre Hamburg (UCC), Hamburg, Germany

<sup>27</sup> Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, WA, USA

<sup>28</sup> Department of Epidemiology, University of Washington, Seattle, WA, USA

<sup>29</sup> Centre for Cognitive Ageing and Cognitive Epidemiology, University of Edinburgh, Edinburgh EH8 9JZ, UK

<sup>30</sup> Psychology, University of Edinburgh, Edinburgh EH8 9JZ, UK

<sup>31</sup> Department of Medicine, Stanford University, Stanford, CA, USA

<sup>32</sup> Department of Cardiology, University Medical Center Groningen, University of Groningen, Groningen, The Netherlands

<sup>33</sup> Department of Complex Trait Genetics, Center for Neurogenomics and Cognitive Research, Amsterdam Neuroscience, VU University Amsterdam, Amsterdam, Netherlands

<sup>34</sup> Institute of Cardiovascular and Medical Sciences, College of Medical, Veterinary and Life Sciences, University of Glasgow, Glasgow, UK

<sup>35</sup> Centre for Cancer Genetic Epidemiology, Department of Oncology, Cambridge Centre, University of Cambridge, Cambridge CB1 8RN, UK

<sup>36</sup> Department of Epidemiology and Biostatistics, Imperial College London, London, UK

<sup>37</sup> MRC-PHE Centre for Environment and Health, Imperial College London, London W2 1PG, UK

<sup>38</sup> National Institute for Health Research Imperial Biomedical

- Research Centre, Imperial College Healthcare NHS Trust and Imperial College London, London, UK
- 39 UK Dementia Research Institute (UK DRI) at Imperial College London, London, UK
- 40 Department of Hygiene and Epidemiology, University of Ioannina Medical School, Ioannina, Greece
- 41 Department of Medical and Molecular Genetics, Indiana University School of Medicine, Indianapolis, IN, USA
- 42 Department of Neuroscience, Icahn School of Medicine at Mount Sinai, New York, NY, USA
- 43 Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, WA, USA
- 44 Department of Psychiatry and Psychotherapy, University Medicine Greifswald, 17475 Greifswald, Germany
- 45 Department of Preventive Medicine, Keck School of Medicine, University of Southern California, Los Angeles, CA, United States
- 46 Department of Public Health and Clinical Medicine, Nutritional research, Umeå University, Umeå, Sweden
- 47 Centre for Genomic and Experimental Medicine, University of Edinburgh, Edinburgh EH4 2XU, UK
- 48 Genetics of Complex Traits, University of Exeter Medical School, Exeter, UK
- 49 University of Southern California, California, CA, USA
- 50 William Harvey Research Institute, Barts and The London School of Medicine and Dentistry, Queen Mary University of London, London EC1M 6BQ, UK
- 51 Centre for Genomic Health, Queen Mary University of London, London EC1M 6BQ, UK
- 52 Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Helsinki, Finland
- 53 Department of Public Health, University of Helsinki, Helsinki, Finland
- 54 Department of Epidemiology, School of Public Health, University of Michigan, Ann Arbor, MI, USA
- 55 Oxford Centre for Diabetes, Endocrinology and Metabolism, University of Oxford, Oxford, UK
- 56 Oxford National Institute for Health Research, Biomedical Research Centre, Churchill Hospital, Oxford, UK
- 57 Department of Public Health Solutions, National Institute for Health and Welfare, FI-00271 Helsinki, Finland
- 58 National Heart and Lung Institute, Imperial College London, London W12 0NN, UK
- 59 Department of Biostatistics, University of Washington School of Medicine, Seattle, WA, USA
- 60 University of Eastern Finland, Finland, Finland
- 61 MRC Epidemiology Unit, Institute of Metabolic Science, University of Cambridge School of Clinical Medicine, Cambridge CB2 0QQ, UK
- 62 Department of Medical Microbiology, Immunology and Cell Biology, Southern Illinois University School of Medicine, Springfield, IL, USA
- 63 Montreal Heart Institute, Montreal, Quebec H1T 1C8, Canada
- 64 Department of Medicine, Faculty of Medicine, Université de Montréal, Montreal, Quebec H3T 1J4, Canada
- 65 NIHR Biomedical Research Centre at Guy's and St Thomas' Foundation Trust, London SE1 9RT, UK
- 66 MRC Human Genetics Unit, MRC Institute of Genetics and Molecular Medicine, University of Edinburgh, Edinburgh, UK
- 67 Queensland Institute for Medical Research, Brisbane, Australia
- 68 Department of Electron Microscopy/Molecular Pathology, The Cyprus Institute of Neurology and Genetics, 1683 Nicosia, Cyprus
- 69 Estonian Genome Center, University of Tartu, Tartu, Estonia
- 70 Institute of Health Informatics, University College London, London, UK
- 71 Department of Epidemiology and Public health, University Hospital of Strasbourg, Strasbourg, France
- 72 Institute of Genetic Epidemiology, Helmholtz Zentrum München - German Research Center for Environmental Health, Neuherberg, Germany
- 73 Department of Medicine I, Ludwig-Maximilians-University Munich, Munich, Germany
- 74 DZHK (German Centre for Cardiovascular Research), Partner Site Munich Heart Alliance, Munich, Germany
- 75 Initiative for Noncommunicable Diseases, Health Systems and Population Studies Division, International Centre for Diarrhoeal Disease Research, Bangladesh (icddr), International Centre for Diarrhoeal Disease Research, Dhaka, Bangladesh
- 76 Institute of Clinical Chemistry and Laboratory Medicine, University Medicine Greifswald, 17475 Greifswald, Germany
- 77 DZHK (German Centre for Cardiovascular Research), Partner Site Greifswald, University Medicine, Greifswald, Germany
- 78 Department of Clinical Biochemistry Herlev Hospital, Copenhagen University Hospital, Herlev Ringvej 74, DK-2730 Herlev, Denmark
- 79 Department of Epidemiology, University of North Carolina, Chapel Hill, NC, USA
- 80 Survey Research Center, Institute for Social Research, University of Michigan, Ann Arbor, MI, USA
- 81 Istituto di Ricerca Genetica e Biomedica, Consiglio Nazionale delle Ricerche (CNR), Monserrato, Cagliari, Italy
- 82 Department of Clinical Genetics, VU University Medical Centre Amsterdam, Amsterdam Neuroscience, Amsterdam, Netherlands
- 83 International Centre for Circulatory Health, Imperial College London, London, UK
- 84 Centre for Non-Communicable Diseases, Karachi, Pakistan
- 85 Genetic and Molecular Epidemiology Unit, Lund University

- Diabetes Centre, Department of Clinical Sciences, Skåne University Hospital, Lund University, SE-214 28, Malmö, Sweden
- <sup>86</sup> Department of Biobank Research, Umeå University, SE-901 87, Umeå, Sweden
- <sup>87</sup> Departments of Psychiatry and Mathematics, Washington University St. Louis, St. Louis, MO, USA
- <sup>88</sup> University of North Carolina, Chapel Hill, NC, USA
- <sup>89</sup> Department of Public Health & Clinical Medicine, Section for Family Medicine, Umeå universitet, SE 90185 Umeå, Sweden
- <sup>90</sup> Department of Cardiovascular Sciences, University of Leicester, Cardiovascular Research Centre, Glenfield Hospital, Leicester LE3 9QP, UK
- <sup>91</sup> National Institute on Aging, National Institutes of Health, Bethesda, MD, USA
- <sup>92</sup> Department of Psychology, University of Amsterdam & Amsterdam Brain and Cognition, University of Amsterdam, Amsterdam, Netherlands
- <sup>93</sup> Alzheimer Scotland Research Centre, University of Edinburgh, Edinburgh EH8 9JZ, UK
- <sup>94</sup> Department of Haematology, University of Cambridge, Cambridge CB2 0PT, UK
- <sup>95</sup> Department of Preventative Medicine, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA
- <sup>96</sup> Wellcome Trust Sanger Institute, Hinxton, Cambridge CB10 1SA, UK
- <sup>97</sup> Department of Medicine, Vanderbilt University, Nashville, TN, USA
- <sup>98</sup> Department of Cardiology, Division Heart and Lungs, University Medical Center Utrecht, Utrecht University, 3508GA Utrecht, The Netherlands
- <sup>99</sup> Department of gerontology and geriatrics, Leiden University Medical Center, Leiden, The Netherlands
- <sup>100</sup> Department of cardiology, Leiden University Medical Center, Leiden, The Netherlands
- <sup>101</sup> Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, 301 Binney Street, Cambridge, MA 02142, USA
- <sup>102</sup> Institute for Community Medicine, University Medicine Greifswald, 17475 Greifswald, Germany
- <sup>103</sup> Clinical Pharmacology, William Harvey Research Institute, Queen Mary University of London, London EC1M 6BQ, UK
- <sup>104</sup> NIHR Barts Cardiovascular Biomedical Research Centre, Barts and The London School of Medicine and Dentistry, Queen Mary University of London, London EC1M 6BQ, UK
- <sup>105</sup> Interfaculty Institute for Genetics and Functional Genomics, University Medicine and Ernst-Moritz-Arndt-University Greifswald, 17475 Greifswald, Germany
- <sup>106</sup> Department of Neuroscience, Psychology and Behaviour, University of Leicester, Leicester, UK
- <sup>107</sup> Department of Biomedical Engineering, The Pennsylvania State University, University Park, Pennsylvania, PA 16802, USA
- <sup>108</sup> Institute of Personalized Medicine, Penn State College of Medicine, Hershey, PA, USA
- <sup>109</sup> Department of Clinical Science, Center for Genetics of Host Defense, University of Texas Southwestern, Dallas, TX, USA
- <sup>110</sup> Department of Cardiology, Ealing Hospital, London North West Healthcare NHS Trust, Middlesex UB1 3HW, UK
- <sup>111</sup> Department of Biostatistics and Epidemiology, University of Pennsylvania, Pennsylvania, PA, USA
- <sup>112</sup> School of Medicine, University of Dundee, Dundee, UK
- <sup>113</sup> Department of Epidemiology and Public Health, Institut Pasteur de Lille, Lille, France
- <sup>114</sup> Department of General and Interventional Cardiology, University Heart Center Hamburg, Hamburg, Germany
- <sup>115</sup> University Medical Center Hamburg Eppendorf, Hamburg, Germany
- <sup>116</sup> William Harvey Research Institute, Barts and The London School of Medicine and Dentistry, Queen Mary University of London, London EC1M 6BQ, UK
- <sup>117</sup> Princess Al-Jawhara Al-Brahim Centre of Excellence in Research of Hereditary Disorders (PACER-HD), King Abdulaziz University, Jeddah 21589, Saudi Arabia
- <sup>118</sup> National Institute for Health Research Blood and Transplant Research Unit in Donor Health and Genomics, Department of Public Health and Primary Care, University of Cambridge, Cambridge CB1 8RN, UK
- <sup>119</sup> EPIMED Research Centre, Department of Medicine and Surgery, University of Insubria at Varese, Varese, Italy
- <sup>120</sup> Department of Epidemiology, UMR 1027- INSERM, Toulouse University-CHU Toulouse, Toulouse, France
- <sup>121</sup> Department of Nutrition, Harvard T. H. Chan School of Public Health, Boston, MA 02115, USA
- <sup>122</sup> Division of Respiratory Medicine and NIHR Nottingham Biomedical Research Centre, University of Nottingham, Nottingham, UK
- <sup>123</sup> Department of Public Health and Clinical Medicine, Skellefteå Research Unit, Umeå University, Umeå, Sweden
- <sup>124</sup> Department of Cardiology, Leiden University Medical Center, Leiden, The Netherlands
- <sup>125</sup> The Interuniversity Cardiology Institute of the Netherlands, Utrecht, The Netherlands
- <sup>126</sup> UKCRC Centre of Excellence for Public Health, Queens, University, Belfast, Belfast, UK
- <sup>127</sup> Medical Research Institute, University of Dundee, Ninewells Hospital and Medical School, Dundee, UK
- <sup>128</sup> Institute of Cardiovascular and Medical Sciences, University of Glasgow, Glasgow, UK
- <sup>129</sup> Department of Twin Research and Genetic Epidemiology, Kings

College London, London SE1 7EH, UK

<sup>130</sup> Population Health Research Institute, St George's, University of London, London SW17 0RE, UK

<sup>131</sup> Department of Genetics, University of Groningen, University Medical Center Groningen, Groningen, The Netherlands

<sup>132</sup> Department of Biostatistics and Epidemiology, Perelman School of Medicine, University of Pennsylvania, Pennsylvania, PA, USA

<sup>133</sup> Center for Non-Communicable Diseases, Karachi, Pakistan

<sup>134</sup> National Institute for Health Research Leicester Respiratory Biomedical Research Centre, Glenfield Hospital, Leicester, UK